

# 線型回帰モデルにおける異常値 の検出と頑健推定によよぼす独立 変数観測値の影響

原田桂一郎

## 目 次

- 1 はじめに
- 2 回帰係数の最小二乗推定に与える観測値の影響
- 3 Hat 行列とその性質
- 4 異常値の検出と独立変数観測値
- 5 回帰係数の頑健推定と独立変数観測値
- 6 むすび
- 付 錄

## 1 はじめに

経済現象を計量経済学的に分析するためには、経済現象を経済変数間の因果関係として捉え、これを明示的に表現したモデル、そして統計データを必要とする。モデルは、現象を構成する本質的な要因を変数とし、理論に基づいて変数の一次式によって現象を表現したものである。統計データは無数の要因から成る現象を観測することによって与えられるものである。

モデルは多くの要因を変数として組み込んで複雑化すれば、現実（データ）と合致したものとすることができます。つまり、理論とデータが整合的であるようなモデルをいくらでも構築することが可能なのである。したがって、計量モデルを用いて理論的命題（仮説）を経験（観測結果）によっ

線型回帰モデルにおける異常値の検出と頑健推定によよばず独立変数観測値の影響で検証することは不可能となる。

計量経済分析において用いられるモデルは、理論を統計データによって検証する、あるいは理論が統計データと整合的であるか否かを調べるためのものではない。経済理論を正しいものと考え、それに基づいて経済関係を計量モデル化する。そして統計データによって推定した経済関係（構造）によって経済現象の数量分析や予測が行われる。経済理論を統計データよりも重視し優先させる。

ここで、現象の観測によって理論が検証できないとすると、経済理論の有意性の規準はどこに求められるか、という重要な問題が存在するが、この点については本稿の主題から離れたものとなるので取り扱わないことにしたい。

さて、本稿で取り扱う計量モデル（單一方程式モデルを考える）は次の型式となる。すなわち、経済変数間の関係を一次式をもって表わした理論モデルは現実の経済現象の近似表現であるから、現実と理論モデルの乖離は誤差項で埋め合はされ線型回帰モデルの型式となる。誤差項は無数の独立な要因（誤差）の和であるとして、中心極限定理により近似的に正規分布にしたがう確率変数と仮定される。正規分布は裾野が薄い分布であるから、大きな誤差、あるいは経済関係を大きく左右するような攪乱要因が発生する確率が低いことを仮定している。さらに、データ観測期間を通じて構造変化<sup>①</sup>が生じないことを仮定すると、最小二乗推定される回帰係数（経済関係）は極めて安定的なものになる。

しかし、データ観測期間中に構造変化がない場合でも、通常でない事態の発生に関連し攪乱要因を含んだ（経済変数の）観測値、および様々な原因による誤差を含んだ観測値の存在は、それ等を推定に使用するか否かによって結果が大いに異なることが、実証分析の経験から明らかとなっていく。こうした観測値は、従属変数の観測値の中に異常値（outliers）となつて出現し、安定した経済関係の推定を妨げる。

このような攪乱要因や誤差が変数から分離できたものとして（現実には

線型回帰モデルにおける異常値の検出と頑健推定によばず独立変数観測値の影響不可能), 回帰モデルの誤差項において処理すると次のようなものとなる。その第一は、少数であるが低い確率で発生する作用の大きな攪乱要因を考慮する、第二は誤差項を裾野の厚い分布(正規分布ではなく、大きな誤差あるいは攪乱要因が生じる確率が低くない)を仮定する、そして第三は両者を合併したものとする、が考えられる。いずれのケースで処理しても、回帰係数の最小二乗推定量の効率は著しく低下する。

異常値が存在する場合の現実的対応は、一定の基準により異常値を検出・除去し、正規母集団から抽出した標本観測値のみを用いて回帰係数の最小二乗推定を行うことである。また、異常値の検出・除去をしないが、最小二乗推定とは別的方式ですすめられる頑健推定(robust estimation)が開発されており、この方法を採用することもできる。

ところが、独立変数の観測値の中に、他の観測値に比べて特に大きな値を示すもの(high leverage)が存在すると、これが異常値の検出と回帰係数の頑健推定に著しい影響を与える。経済関係を線型回帰モデルによって分析する場合、独立変数観測値のhigh leverageは、究極的に安定した経済関係(回帰係数)を推定することを妨げる。

以下において、線型回帰モデルの係数の最小二乗推定に与える観測値の影響、独立変数観測値のhigh leverageの検出、異常値の検出と頑健推定におけるhigh leverageの影響、について考察する。これらを通して、観測値から異常値やhigh leverageを検出・除去して、安定した回帰係数を推定する方途を示す。

## 注

- ① 構造変化の問題については、Harada, K. and R. Shirogane : "Recursive Estimation Analysis for Examining the Stability of Regression Relationships," *Memoirs of the Kokushikan University Computer Center*, No. 6 (1985), 1-20.

## 2 回帰係数の最小二乗推定に与える観測値の影響

次の標準線型回帰モデルを考える。

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \quad (2-1)$$

$\mathbf{y}$ ：従属（被説明）変数ベクトル  $\mathbf{y} = (y_1 \cdots y_n)'$ ,

$\mathbf{X}$ ：独立（説明）変数行列 ( $n \times p$ ),  $\text{rank}(\mathbf{X}) = p$ ,  $(2-2)$

$\boldsymbol{\beta}$ ：回帰係数ベクトル  $\boldsymbol{\beta} = (\beta_1 \cdots \beta_p)'$ ,

$\mathbf{u}$ ：誤差ベクトル  $\mathbf{u} = (u_1 \cdots u_n)'$ ,

$$E(\mathbf{u}) = 0, \quad V(\mathbf{u}) = \sigma^2 \mathbf{I}, \quad \mathbf{I}(n \times n),$$

そして、(2-1) が推定されたものを、

$$\mathbf{y} = \mathbf{X}\mathbf{b} + \mathbf{e}, \quad (2-3)$$

$\mathbf{y}$  : (2-1) に同じ,

$\mathbf{X}$  : (2-1) に同じ,

$\mathbf{b}$  :  $\boldsymbol{\beta}$  の最小二乗推定量,

$$\mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}, \quad (2-4)$$

$\mathbf{e}$  : 残差ベクトル  $\mathbf{e} = (e_1 \cdots e_n)'$

$S^2$  : 誤差分散  $\sigma^2$  の推定量,

$$S^2 = \mathbf{e}'\mathbf{e}/(n-p), \quad (2-5)$$

と表わす。さらに次の表記を与えておく。

$\mathbf{x}_i$  :  $\mathbf{X}$  行列の第  $i$  行 (ベクトル)

$$\mathbf{X} = [\mathbf{x}_1 \mathbf{x}_2 \cdots \mathbf{x}_n]', \quad \mathbf{X}' = [\mathbf{x}_1' \mathbf{x}_2' \cdots \mathbf{x}_n'],$$

$$\mathbf{X}'\mathbf{X} = \sum_{i=1}^n \mathbf{x}_i' \mathbf{x}_i$$

$\mathbf{X}(i) = \mathbf{X}$  行列から第  $i$  行を除いた行列,

$\mathbf{y}(i)$  : ベクトルから第  $i$  要素を除いたベクトル,

線型回帰モデルにおける異常値の検出と頑健推定によばず独立変数観測値の影響

$\mathbf{b}(i) : \mathbf{X}$  と  $\mathbf{y}$  から第  $i$  行, 第  $i$  要素を除いた場合の  $\beta$  の最小二乗推定量,

$$\mathbf{b}(i) = [\mathbf{X}'(i)\mathbf{X}(i)]^{-1}\mathbf{X}'(i)\mathbf{y}(i), \quad (2-6)$$

$S^2(i) : \mathbf{X}$  と  $\mathbf{y}$  から第  $i$  行, 第  $i$  要素を除いた場合の  $\sigma^2$  の推定量,

$$S^2(i) = [\mathbf{y}(i) - \mathbf{X}(i)\mathbf{b}(i)]'[\mathbf{y}(i) - \mathbf{X}(i)\mathbf{b}(i)] / (n-p-1). \quad (2-7)$$

ここで, 回帰係数の最小二乗推定量にデータ  $\{(y_i, \mathbf{x}_i)\}_{i=1}^n$  が与える影響を調べるため, 第  $i$  観測値にウェイトを付けたときの最小二乗推定量  $\mathbf{b}(w_i)$  を考察する。

$$\mathbf{b}(w_i) = (\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}\mathbf{y}, \quad (2-8)$$

$$\mathbf{W} = \begin{pmatrix} 1 & 0 & \cdots & \cdots & 0 \\ 0 & \cdot & & & \vdots \\ \vdots & & 1 & & \vdots \\ \vdots & & w_i & & \vdots \\ \vdots & & & 1 & \vdots \\ \vdots & & & & \vdots \\ 0 & \cdots & \cdots & \cdots & 1 \end{pmatrix}$$

と与えられ, これを  $w_i$  について微分すると,

$$\frac{\partial \mathbf{b}(w_i)}{\partial w_i} = \frac{(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i'(y_i - \mathbf{x}_i\mathbf{b})}{[1 - (1-w_i)\mathbf{x}_i'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i']^2} \quad (2-9)$$

さらに,  $w_i=1$  とすると,

$$\left. \frac{\partial \mathbf{b}(w_i)}{\partial w_i} \right|_{w_i=1} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i'(y_i - \mathbf{x}_i\mathbf{b}) \quad (2-10)$$

を得る。この式により  $(y_i, \mathbf{x}_i)$  が回帰係数の最小二乗推定量に与える影響を調べることができる<sup>②</sup>。(2-10) の両辺の左側に  $\mathbf{x}_i$  を乗ずると,

$$\mathbf{x}_i \cdot \frac{\partial \mathbf{b}(w_i)}{\partial w_i} \Big|_{w_i=1} = \mathbf{x}_i (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_i' (y_i - \mathbf{x}_i b) \quad (2-11)$$

のような形になる。

$\mathbf{x}_i (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_i'$  は行列  $\mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}'$  の第  $i$  対角要素であり、 独立変数データの情報を与え、 残差  $e_i = y_i - \mathbf{x}_i b$  は従属変数データの情報を与える。 $\mathbf{x}_i (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_i'$  の値が大であるときの独立変数観測値は high leverage,  $e_i$  が大であるときの従属変数観測値は異常値と呼ばれる<sup>③</sup>。

いま、  $L(u_i | \mathbf{x}_i) = N(0, \sigma^2)$  を仮定すると、 データ  $\{(y_i, \mathbf{x}_i)\}_{i=1}^n$  は多変量正規分布  $\mathbf{y} \sim N(\mathbf{X}\beta, \sigma^2 \mathbf{I})$  から独立に抽出されたものである。このときの回帰係数の最小二乗推定量の効率は大である。しかし、 観測値に異常値や high leverage が存在すると、その効率は低下する<sup>④</sup>。現実的には、推定値  $\mathbf{b}$  が歪んだ値 ( $\beta$  からかけ離れた値)となってしまう。したがって、4節に記述する所定の方法に基づいて異常値や high leverage 独立変数観測値を検出し、それ等を除いてから回帰係数の最小二乗推定を行うことが考えられる。

## 注

② (2-10) は Belsley et al. [3] により提示された回帰係数の最小二乗推定量の influence function である。

③ 母集団分布とは異なる別の母集団から生じた観測値を異常値と名づける。

④ 効率の低下についての詳細な説明は、Fisher 情報行列を用いて行うことができる。

## 3 Hat 行列とその性質

異常値や独立変数の high leverage が回帰係数の最小二乗推定量を歪めたものとする。本節では、独立変数が従属変数の予測値および残差に影響することを示し、独立変数による行列  $\mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}'$  の対角要素の性質

線型回帰モデルにおける異常値の検出と頑健推定によばず独立変数観測値の影響を検討する。

まず、従属変数の予測値は  $\hat{\mathbf{y}} = \mathbf{X}\mathbf{b}$  で与えられるから、

$$\hat{\mathbf{y}} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}. \quad (3-1)$$

ここで、 $\hat{\mathbf{y}}$  の要素  $\hat{y}_i$  が一次関数となることから、(3-1) を

$$\begin{aligned} \hat{\mathbf{y}} &= \mathbf{H}\mathbf{y}, \\ \mathbf{H} &= \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}', \end{aligned} \quad (3-2)$$

のように表わす。行列  $\mathbf{H}(n \times n)$  は  $\mathbf{y}$  により  $\hat{\mathbf{y}}$  の値を定める関数(写像)の作用をなし、これを“Hat 行列”と称す。

幾何学的にみると、ベクトル  $\mathbf{y}$  と  $\mathbf{X}$  の列を  $n$  次元ユーベリット空間の点で表わすなら、点  $\mathbf{X}\beta$ (列ベクトル  $\mathbf{x}_j$  の一次結合  $\sum_{j=1}^p \beta_j \mathbf{x}_j$  である) は  $n$  次元部分空間をなす。ベクトル  $\hat{\mathbf{y}}$  は  $\mathbf{y}$  に最も近いその部分空間の点であり、それはまた  $\mathbf{y}$  のその部分空間への直交射影(perpendicular projection)である。したがって  $\mathbf{H}$  は射影行列(projection matrix)である(Hoaglin and Welsch [12])。

また、行列  $\mathbf{H}$  は残差ベクトルにも現われて、

$$\begin{aligned} \mathbf{e} &= \mathbf{y} - \mathbf{X}\mathbf{b} \\ &= [\mathbf{I} - \mathbf{H}]\mathbf{y}, \end{aligned} \quad (3-3)$$

となる。独立変数から成る行列  $\mathbf{H}$  およびその要素  $h_{ij}$  のもつ性質は、Belsley et al. [3] と Hoaglin and Welsch [12] の導出した結果に若干の補追を加えると以下のようである。

まず、 $\hat{\mathbf{y}}$  の要素  $\hat{y}_i$  は、

$$\begin{aligned}
 \hat{y}_i &= \mathbf{x}_i \mathbf{b} = \mathbf{x}_i (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y} \\
 &= h_{ii} y_i + \sum_{j \neq i} h_{ij} y_j, \\
 h_{ii} &\equiv \mathbf{x}_i (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_i', \quad (\text{行列 } \mathbf{H} \text{ の第 } i \text{ 対角要素}) \\
 h_{ij} &\equiv \mathbf{x}_i (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_j', \quad (\text{行列 } \mathbf{H} \text{ の } ij \text{ 要素})
 \end{aligned} \tag{3-4}$$

のように表わされる。(3-4) は、行列  $\mathbf{H}$  の要素  $h_{ij}$  を梃子に  $y_j$  が  $\hat{y}_i$  に影響を及ぼすと解される。 $\hat{y}_i$  に及ぼす  $y_i$  の直接的な影響は対角要素  $h_{ii}$  を梃子にして現われる。

独立変数の観測値を採取していくと、行列  $\mathbf{X}$  のどの行ベクトルが、 $\hat{\mathbf{y}}$  の該当する要素の値を大きくしているかを判定できる。行列  $\mathbf{H}$  が、その “high leverage point,  $i$ ” を識別するために採用される。特に、その対角要素  $h_{ii}$  がそのためには有効である。次に  $h_{ii}$  の性質を検討する。

行列  $\mathbf{H}$  は対称かつ巾等 (symmetric and idempotent) であるから ( $\mathbf{H}^2 = \mathbf{H}$ )、対角要素は

$$h_{ii} = \sum_{j=1}^n h_{ij}^2 = h_{ii}^2 + \sum_{j \neq i} h_{ij}^2, \tag{3-5}$$

となり、また  $0 \leq h_{ii} \leq 1$  である。巾等行列の固有値 (eigenvalues) は 0 か 1 のいずれかであり、また零ではない固有値の数はその行列の階数 (rank) に等しいから、行列  $\mathbf{H}$  の場合には、 $\text{rank}(\mathbf{H}) = \text{rank}(\mathbf{X}) = p$ 、である。巾等行列の rank と trace は等しいから、 $\text{trace}(\mathbf{H}) = p$ 、すなわち、

$$\sum_{i=1}^n h_{ii} = p \tag{3-6}$$

なお、行列  $\mathbf{H}$  の対角要素の大きさの平均は、 $p/n$ 、となる。

ところで、 $\mathbf{l} = [1 \cdots 1]', \mathbf{U}' = \mathbf{I}, \bar{\mathbf{y}} = \frac{1}{n} \mathbf{l}' \mathbf{y}$ 、として  $(\hat{\mathbf{y}} - \bar{\mathbf{y}})$  を展開す

線型回帰モデルにおける異常値の検出と頑健推定によばず独立変数観測値の影響ると、

$$\hat{\mathbf{y}} - \bar{\mathbf{y}} = \mathbf{H}\mathbf{y} - l\bar{\mathbf{y}} \\ = \left[ \mathbf{H} - \frac{1}{n} \mathbf{I} \right] \mathbf{y} = \tilde{\mathbf{H}}\mathbf{y} \quad (3-7)$$

$$\tilde{\mathbf{H}} = \mathbf{H} - \frac{1}{n} \mathbf{I} \quad (3-8)$$

となる。行列  $\tilde{\mathbf{H}}$  の要素  $\tilde{h}_{ij}$  は (3-8) から、

$$\tilde{h}_{ij} = h_{ij} - \frac{1}{n} \quad (3-9)$$

のようになる。ここで行列  $\mathbf{X}$  の各列  $\mathbf{x}_j = (x_{1j} \dots x_{nj})'$  について平均をとり、 $\bar{\mathbf{x}}_j = \frac{1}{n} \mathbf{l}' \mathbf{x}_j$ 、ベクトル  $\bar{\mathbf{x}} = (\bar{x}_1 \bar{x}_2 \dots \bar{x}_p)$  をつくる。そして、

$$\tilde{\mathbf{X}} = \begin{pmatrix} \mathbf{x}_1 - \bar{\mathbf{x}} \\ \mathbf{x}_2 - \bar{\mathbf{x}} \\ \vdots \\ \mathbf{x}_n - \bar{\mathbf{x}} \end{pmatrix} = \begin{pmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \cdots & x_{1p} - \bar{x}_p \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & \cdots & x_{2p} - \bar{x}_p \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} - \bar{x}_1 & x_{n2} - \bar{x}_2 & \cdots & x_{np} - \bar{x}_p \end{pmatrix} = \begin{pmatrix} \tilde{\mathbf{x}}_1 \\ \tilde{\mathbf{x}}_2 \\ \vdots \\ \tilde{\mathbf{x}}_n \end{pmatrix} \quad (3-10)$$

$$\tilde{\mathbf{H}} = \tilde{\mathbf{X}} (\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}', \quad (3-11)$$

を定義する。行列  $\tilde{\mathbf{H}}$  の対角要素  $\tilde{h}_{ii}$  は、

$$\tilde{h}_{ii} = (\mathbf{x}_i - \bar{\mathbf{x}}) (\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1} (\mathbf{x}_i - \bar{\mathbf{x}})' \quad (3-12)$$

であるから、行列  $\mathbf{H}$  の対角要素は (3-9) から

$$h_{ii} = \frac{1}{n} + (\mathbf{x}_i - \bar{\mathbf{x}})(\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1}(\mathbf{x}_i - \bar{\mathbf{x}})' \quad (3-13)$$

のように与えられる。

$h_{ii}$  は  $(\mathbf{x}_i - \bar{\mathbf{x}})$  が大であるとき、即ち独立変数観測値の中でその平均値よりもはるかに大きな値を取るととき、大きくなる（1に近づく）。また  $(\mathbf{x}_i - \bar{\mathbf{x}})$  が小であるとき、小さな値となる（零に近づく）。

#### 4 異常値の検出と独立変数観測値

線型回帰モデルでは誤差項を確率変数と仮定するので、(2-1)により従属変数は確率変数となる（独立変数は確率変数ではない）。従属変数の観測値が、（誤差項について）仮定した確率分布と異なる別の分布の母集団から抽出された観測値（異常値）か否かの検定を行うとき、誤差項に正規性の仮定を置けば、Studentized residuals を用いることができる。この Studentized residuals は、

$$e_{si} = \frac{e_i}{S(1-h_{ii})^{\frac{1}{2}}} \quad (4-1)$$

と与えられるが、分子の残差  $e_i$  にも行列  $\mathbf{H}$  の対角要素  $h_{ii}$  が存在し、異常値の検出に影響を及ぼしている。

残差は (2-3) と (3-4) を用いて、

$$\begin{aligned} e_i &= y_i - \hat{y}_i \\ &= (1-h_{ii})y_i - \sum_{j \neq i} h_{ij}y_j, \end{aligned} \quad (4-2)$$

と表わすことができる。しかし、異常値  $y_i$  は常に該当する残差  $e_i$  を大とするとはかぎらない。 $h_{ii}$  が零に近い小さな値であるとき  $e_i$  は大とな

線型回帰モデルにおける異常値の検出と頑健推定によばず独立変数観測値の影響するが、 $h_{ii}$  が 1 に近い大きな値であるときには  $e_i$  は小さな値となる。また、異常値  $y_i$  は該当する残差  $e_i$  ではなく、他の残差を大とする可能性がある。すなわち、もし  $h_{mi}$  の値が大きければ、 $y_i$  により残差  $y_m - \hat{y}_m$  の値が大きくなるからである。

残差  $e_i$  は  $h_{ii}$  の減少関数となることから、Studentized residuals  $e_{si}$  による異常値検定の検出力は、 $h_{ii}$  の減少関数である<sup>⑤</sup>。 $h_{ii}$  が大である場合の異常値を検出する能力は、 $h_{ii}$  が小である場合に比べて劣ることになる。

ところで、線型回帰モデルの回帰係数の最小二乗推定量の漸近特性（大標本特性）については、Anderson [1, pp.23-27] と Eickel [8] によって取り扱かわれ、その漸近的正規性が証明されているが、Huber [15] は線型回帰係数の最小二乗推定量が漸近的正規性を保持するための必要十分条件は

$$\lim_{n \rightarrow \infty} (\max_i h_{ii}) = 0, \quad i=1, 2, \dots, n \quad (4-3)$$

であることを示している。 $h_{ii}$  の値が小さいことは、線型回帰係数の最小二乗推定量の漸近的正規性のためにも必要なことである。さらに Huber は (4-3) の条件はまた ‘一種の “design (行列  $\mathbf{X}$ ) の頑健性” を保証する’ とし、また ‘この design の頑健性はすべての  $i$  について  $h_{ii} = p/n$  とすることによって最大化され、そのときに、design 行列  $\mathbf{X}$  が均齊している’ としている (Huber [15, p.804])。また、Box and Draper [5] は異常値に対して design を鋭敏でなくするためには、 $h_{ii}$  を小さくすることだとしている。これらのこととは、(3-4) の関係から、異常値が存在するとき  $h_{ii}$  が大であると、これが梃子の作用をして予測値  $\hat{y}_i$  を大きくしてしまうことを警戒したものである。

(3-13) から明らかなように、 $h_{ii}$  が大であることは、独立変数の第  $i$  観測値  $(x_{i1}, x_{i2}, \dots, x_{ip})$  が大きい (high leverage) ことを示している。独

線型回帰モデルにおける異常値の検出と頑健推定によばず独立変数観測値の影響立変数の high leverage 観測値は Studentized residuals による異常値検出の能力を損い、また従属変数の予測値を肥大化する。経済変数のデータは時系列で与えられることが一般的であり、繰り返し測定したりコントロールすることが不可能であるから、行列  $X$  は固定されたものとなる。したがって  $h_{ii}$  を計算して、その値が大きい場合には該当する観測値  $(x_{i1}, \dots, x_{ip})$  を除外しておくことが、上記の事態を回避するために適切な措置といえよう。

そのためには、 $h_{ii}$  がどの程度の大きさとなつたときに独立変数の観測値を削除するかを判定するための基準を設定する必要がある。Belsley et al. [3] は次のような方法で  $h_{ii} > 2p/n$  の場合に  $i$  番目の独立変数データを high leverage だと定めている。

独立変数  $x_i$  が  $\nu$  次元の相互独立な正規分布にしたがうと仮定すると、 $h_{ii}$  について正確な分布を計算することができる。独立変数についてそうした仮定を置くことは、現実的妥当性はないが、上記の目的のみに限定して導出した結果を用いる。その仮定の下で  $h_{ii}$  について、

$$F = \frac{(n-p)}{(p-1)} \frac{\{h_{ii} - \left(\frac{1}{n}\right)\}}{(1-h_{ii})} \quad (4-4)$$

の関係は自由度  $p-1$  および  $n-p$  の  $F$  分布にしたがうことが示される<sup>⑥</sup>。 $p > 10$  および  $n-p > 50$  について、 $F$  値の 95% ポイントが 2 以下であるから、 $2p/n$  ( $h_{ii}$  の平均値の 2 倍) を  $x_i$  のカット・オフ基準とする。 $p$  が多く  $n$  が少なく、 $p/n > 0.4$  の場合には基準が高くなり過ぎるので、すべての独立変数観測値に疑惑が生じ、 $p$  が少ない場合にはカット・オフの対象となる独立変数観測値が少くなってしまう。しかし、 $h_{ii}$  が  $2p/n$  を越えた場合、 $i$  番目の独立変数観測値の位置を leverage point と称す。

**注**

⑥ Cook [6] は mean shift modelにおいて、第  $i$  観測値が異常値ではないという仮説を検討するために、Gentleman and Wilk [9] の導出した統計量を基にF-統計量を算出している。この統計量は Studentized residuals の単調增加関数となっている。よって、この F-統計量による検定の検出力も  $h_{ii}$  の減少関数である。

⑥ 詳細は付録を参照。

## 5 回帰係数の頑健推定と独立変数観測値

前節の方法のように、一定規準を越える異常値を除外して回帰係数を最小二乗推定するのではなく、頑健推定 (robust estimation) と称する方法が開発されている。

現実問題への統計分析の経験を通して、観測値の母集団分布に正規性を仮定することに疑問が生じる場合がある。その中から生み出された頑健推定は、観測値の母集団を正規分布と仮定せず、また残差平方和最小化の規準を探らない推定方法である。その主たる狙いは、非常に大きな観測誤差に対して安全装置を設け、異常値が推定値に及ぼす影響を一定範囲におさめ、大きな異常値を隔離し、なおかつモデル（本稿の場合は線型回帰モデル）を最適なものとすることである。モデルに対する観測値の影響が、ほとんど検出不可能なものから、甚だしく大きなものまであることを考え、大半の観測値の動きに着目する。観測値母集団に特定の分布を仮定せず、また Gauss-Markov 定理を用いて分布の平均値を推定することをもって最適化を図るものでもない。

Huber [15] は回帰係数の頑健推定の方式を次のように示している。残差平方和の最小化、

$$\sum_{i=1}^n (y_i - \mathbf{x}_i \boldsymbol{\beta})^2 = \min. \quad (5-1)$$

線型回帰モデルにおける異常値の検出と頑健推定によよばず独立変数観測値の影響に代って、残差についての特定の関数を最小

$$\sum_{i=1}^n \rho(y_i - \mathbf{x}_i \boldsymbol{\beta}) = \min \quad (5-2)$$

とするような  $\boldsymbol{\beta}$  を推定する。ここで、 $\rho(t)$  は凸関数 (convex function) であり Huber は、

$$\begin{aligned} \rho(t) &= \frac{1}{2} t^2 & |t| < c \\ &= c|t| - \frac{1}{2} c^2 & |t| \geq c \end{aligned} \quad (5-3)$$

を提示している。 $c$  の値は許容しようと定めた  $y_i$  の大きさに依って設定される。推定値は  $c$  の値に対して敏感に反応しないので、誤差の標準偏差の 1 から 2 倍に  $c$  を定めるのが合理的だとしている (Huber [16])。

関数  $\rho(t)$  が微分可能であれば、 $\psi(t) = \rho'(t)$ 、(5-2) を満足するための必要条件は、

$$\sum_{i=1}^n x_{ij} \psi(y_i - \mathbf{x}_i \boldsymbol{\beta}) = 0, \quad j=1, 2, \dots, P \quad (5-4)$$

である。

一般に、(5-4) は非線型方程式となるから、解を得るために反復計算を必要とする。反復計算による推定値の求め方として, Newton 法, Bickel [4] による方法, Beaton and Tukey [2] による re-weighted least squares がある。いずれの方法を探るにしろ、関数  $\rho(t)$  では極小値に収束することのみが保証されるだけであるから、特性の良い初期値の選定が必要となる。

Holland and Welsch [13] は最小絶対値残差推定による初期値を用いて

線型回帰モデルにおける異常値の検出と頑健推定によばず独立変数観測値の影響 re-weighted least squares を行うことを勧めている。最小絶対値残差推定量は漸近的正規性を有すので、反復計算により求めた回帰係数の頑健推定量もこの漸近特性を保持する。

ところで、Hampel [10] は各観測値  $\{(y_i, \mathbf{x}_i)\}$  が頑健推定に与える影響を示す influence function を定義しているが、Krasker and Welsch [18] にしたがってそれを次のように表わす。

回帰係数  $\boldsymbol{\beta}$  とその推定量  $\mathbf{b}$  について  $\mathbf{b} - \boldsymbol{\beta}$  を考える<sup>⑦</sup>。そして、

$$\lim_{n \rightarrow \infty} Pr \left[ \sqrt{n} \left\{ (\mathbf{b} - \boldsymbol{\beta}) - \frac{1}{n} \sum_{i=1}^n \Omega(y_i, \mathbf{x}_i) \right\} \right] = 0 \quad (5-5)$$

を満足する関数  $\Omega(y_i, \mathbf{x}_i)$  は推定量  $\mathbf{b}$  に対して  $i$  番目の観測値が与える影響を示す influence function である。もし、 $E[\Omega(\mathbf{y}, \mathbf{X})] = 0$ ,  $E|\Omega(\mathbf{y}, \mathbf{X})|^2 < \infty$  ならば、(4-5) により推定量は一致漸近的正規性を有し、その漸近分散は

$$V = E[\Omega(\mathbf{y}, \mathbf{X})\Omega(\mathbf{y}, \mathbf{X})'] \quad (5-6)$$

となる。

Huber [15] の提示したタイプの回帰係数の頑健推定量の influence function は

$$\Omega(y_i, \mathbf{x}_i) = \Psi[y_i - \mathbf{x}_i' \boldsymbol{\beta}] (\mathbf{X}' \mathbf{X})^{-1} \mathbf{x}_i' \quad (5-7)$$

のようになる。 $\Psi(o)$  が大きな値の残差を抑えて  $\mathbf{x}_i'$  が乗せられているから、観測値の影響は任意に大となる。回帰係数の頑健推定では、独立変数の観測値中に他よりも特に大きな値のものがあるときには注意を要す。この点については Huber [15] が頑健推定にあたって、独立変数行列に課せられる条件として、

$$\max_{1 \leq i \leq n} h_{ii} = \ll 1 \quad (5-8)$$

を挙げている。

関数  $\rho(t)$  については幾つかのものが提示されており、モンテカルロ実験により各々の関数による頑健推定量の漸近的効率の比較的検討がなされている (Holland and Welsch [12])。どのような関数  $\rho(t)$  をもってしても、回帰係数の頑健推定の influence function は (5-7) の形になるから、独立変数の high leverage の影響は除去できない。

現実的な対応策として、 $h_{ii}$  ( $i=1 \cdots n$ ) を計算し、 $2p/n$  を越える high leverage 独立変数観測値を除いて頑健推定を行うことが考えられる。

### 注

⑦  $b$  が最小二乗推定量であり、また誤差項に正規性  $u \sim N(0, \sigma^2 I)$  の仮定を置くと、

$$\begin{aligned} b &\sim N[\beta, \sigma^2(X'X)^{-1}], \\ b - \beta &= (X'X)^{-1}X'u, \\ E(b - \beta) &= 0, \\ \text{Var}(b - \beta) &= E[(b - \beta)(b - \beta)'] = \sigma^2(X'X)^{-1} \end{aligned}$$

## 5 むすび

線型回帰モデルの回帰係数の最小二乗推定量は一定条件の下では効率の大なる推定量である。しかし、観測値の母集団分布が裾野の厚いもの（正規分布ではない）であったり、観測値に異常値が存在するとき、また両者が合併したときには、最小二乗推定量の効率は著しく低下する。

観測値母集団に正規性が仮定できる場合、Studentized residuals を用い異常値を検出、除去した後に最小二乗推定を行う。正規性の仮定が設定できず、かつ異常値の影響を排除しようとする場合、頑健推定の方法が開発されている。いずれの方法においても、high leverage 独立変数観測値

線型回帰モデルにおける異常値の検出と頑健推定によばす独立変数観測値の影響の影響を除去できないので,  $h_{ii}$  の計算によって high leverage point を検出し, 該当する独立変数観測値を除外することが現実的対応であろう。

回帰係数の推定にあたって, 異常値や high leverage 独立変数の影響を除去するため, 極く一部の観測値が及ぼす影響を制限しながら回帰係数の推定を行い, その推定量が漸近的正規性を保有しつつその漸近分散ができるだけ小さいものとする試みがなされている (Krasker and Welsch [18])。経済変数間の関係を回帰分析する場合, この推定方法を採用する問題については次の課題としたい。

## 付 錄

### 1 $\tilde{\mathbf{x}}_i$ の分布

いま  $\mathbf{X}$  を多次元正規分布にしたがう確率変数と仮定する。そして次のように  $\tilde{\mathbf{X}}$  を定義する。

$$\tilde{\mathbf{X}} = \begin{pmatrix} \mathbf{x}_1 - \bar{\mathbf{x}} \\ \mathbf{x}_2 - \bar{\mathbf{x}} \\ \vdots \\ \mathbf{x}_i - \bar{\mathbf{x}} \\ \vdots \\ \mathbf{x}_n - \bar{\mathbf{x}} \end{pmatrix} = \begin{pmatrix} \tilde{\mathbf{x}}_1 \\ \tilde{\mathbf{x}}_2 \\ \vdots \\ \tilde{\mathbf{x}}_i \\ \vdots \\ \tilde{\mathbf{x}}_n \end{pmatrix} = \begin{pmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \cdots & x_{1p} - \bar{x}_p \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & \cdots & x_{2p} - \bar{x}_p \\ \vdots & \vdots & & \vdots \\ x_{i1} - \bar{x}_1 & x_{i2} - \bar{x}_2 & & x_{ip} - \bar{x}_p \\ \vdots & \vdots & & \vdots \\ x_{n1} - \bar{x}_1 & x_{n2} - \bar{x}_2 & & x_{np} - \bar{x}_p \end{pmatrix}$$

$$\bar{\mathbf{x}} = (\bar{x}_1 \bar{x}_2 \cdots \bar{x}_p),$$

$$\bar{x}_j = \frac{1}{n} l' \mathbf{x}_j, \quad j = 1 \cdots p, \quad l = (1 \cdots 1)' (n \times 1),$$

$$\mathbf{x}_j = (x_{1j} x_{2j} \cdots x_{nj})',$$

$$\tilde{\mathbf{X}}' \tilde{\mathbf{X}} = \sum_{i=1}^n \tilde{\mathbf{x}}_i' \tilde{\mathbf{x}}_i,$$

$$\bar{\tilde{x}} = \frac{1}{n} \sum_{i=1}^n \tilde{x}_i = 0, \quad (\because \sum \tilde{x}_i = 0)$$

$$\bar{\tilde{x}}(i) = \frac{1}{n-1} \sum_{k \neq i} \tilde{x}_k = (n\bar{\tilde{x}} - \tilde{x}_i)/(n-1) = (\bar{\tilde{x}} - \tilde{x}_i/n)/(n-1)/n$$

ここで、 $\tilde{x}_i$  の影響を調べるために、 $\tilde{X}$  から  $\tilde{x}_i$  を除いた  $\tilde{X}(i)$  の分散は、

$$V\{\tilde{X}(i)\} = \frac{1}{n-1} \sum_{k \neq i} \tilde{x}_k' \tilde{x}_k - \bar{\tilde{x}}'(i) \bar{\tilde{x}}(i) = \frac{1}{n-1} \left( \sum_{i=1}^n \tilde{x}_i' \tilde{x}_i - \tilde{x}_i' \tilde{x}_i \right) - \bar{\tilde{x}}'(i) \bar{\tilde{x}}(i)$$

また、 $\tilde{X}$  の分散は  $\bar{\tilde{x}} = 0$  であるから

$$V(\tilde{X}) = \frac{1}{n-P} \sum_{i=1}^n \tilde{x}_i' \tilde{x}_i$$

である。 $V\{\tilde{X}(i)\}$  と  $V(\tilde{X})$  についての Wilk's  $A$  statistics (Rao [19, p. 570]) は、

$$A(\tilde{x}_i) = \frac{\det\{\tilde{X}' \tilde{X} - (n-1) \bar{\tilde{x}}'(i) \bar{\tilde{x}}(i) - \tilde{x}_i' \tilde{x}_i\}}{\det(\tilde{X}' \tilde{X})}$$

となる。分子は

$$\det\left\{\tilde{X}' \tilde{X} - \frac{n^2}{(n-1)} \left(\bar{\tilde{x}} - \frac{\tilde{x}_i}{n}\right)' \left(\bar{\tilde{x}} - \frac{\tilde{x}_i}{n}\right) - \tilde{x}_i' \tilde{x}_i\right\}$$

となり、 $\bar{\tilde{x}} = 0$  であるから

$$\det \left( \tilde{\mathbf{X}}' \tilde{\mathbf{X}} - \frac{n}{n-1} \tilde{\mathbf{x}}_i' \tilde{\mathbf{x}}_i \right)$$

さらに,

$$\begin{aligned} \det \left( \tilde{\mathbf{X}}' \tilde{\mathbf{X}} - \frac{n}{n-1} \tilde{\mathbf{x}}_i' \tilde{\mathbf{x}}_i \right) &= \left( 1 - \frac{n}{n-1} \tilde{\mathbf{x}}_i' (\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{x}}_i \right) \det(\tilde{\mathbf{x}}' \tilde{\mathbf{x}}) \\ &= \left( 1 - \frac{n}{n-1} \tilde{h}_{ii} \right) \det(\tilde{\mathbf{X}}' \tilde{\mathbf{X}}), \end{aligned}$$

よって  $A(\tilde{\mathbf{x}}_i)$  は,

$$A(\tilde{\mathbf{x}}_i) = 1 - \frac{n}{n-1} \tilde{h}_{ii} = \frac{n}{n-1} (1 - h_{ii})$$

となる。

$A(\tilde{\mathbf{x}}_i)$  は  $\chi^2$  分布にしたがう統計量であり、この値が大ならば（小ならば） $\tilde{\mathbf{x}}_i$  の影響が小（ $\tilde{\mathbf{x}}_i$  の影響大）となる。 $\tilde{\mathbf{X}}$  は  $p$  次元正規分布からの  $n$  個の独立な標本で成るから、 $A(\tilde{\mathbf{x}})$  は  $F$  分布にして次のように表わせる(Rao [19, p.570])。

$$\frac{n-p}{p-1} \left[ \frac{1-A(\tilde{\mathbf{x}}_i)}{V(\tilde{\mathbf{x}}_i)} \right] \sim F_{p-1, n-p}$$

したがって、

$$\frac{n-p}{p-1} \frac{[h_{ii} - (1/n)]}{(1-h_{ii})} \sim F_{p-1, n-p}.$$

(1986年3月)

## REFERECE

- [ 1 ] Anderson, T. W. : *The Statistical Analysis of Time Series*. New York : John Wiley & Sons, 1971.
- [ 2 ] Beaton, A. E. and J. W. Tukey. "The Fitting of Power Series, Meaning Polynomials, Illustrated on Band-Spectroscopic Data," *Technometrics*, 16 (1974), 147-185.
- [ 3 ] Belsley, D. A., E. Kuh and R. E. Welsch : *Regression Diagnostics ; Identifying Influential Data and Sources of Collinearity*. New York : John Wiley & Sons, 1980.
- [ 4 ] Bickel, P. J. : "One Step Huber Estimates in the Linear Model," *Journal of the American Statistical Association*, 70 (1975), 428-434.
- [ 5 ] Box, G. E. P. and N. R. Draper : "Robust Designs", *Biometrika*, 62(1975), 347-351.
- [ 6 ] Cook, R. D. : "Influential Observations in Linear Regression," *Journal of the American Statistical Association*, 74 (1979), 169-174.
- [ 7 ] Cook, R. D. and S. Weisberg : "Characterizations of an Empirical Influence Function for Detecting Influential Cases in Regression," *Technometrics*, 22 (1980), 495-508.
- [ 8 ] Eicker, F. : "Asymptotic Normarity and Consistency of the Least Squares Estimators for Families of Linear Regressions", *The Annals of Mathematical Statistics*, 34 (1963), 447-456.
- [ 9 ] Gentleman, J. F. and M. B. Wilk : "Detecting Outliers II. Supplementing the Direct Analysis of Residuals", *Biometrics*, 31 (1975), 387-410.

線型回帰モデルにおける異常値の検出と頑健推定におよぼす独立変数観測値の影響

- [10] Hampel, F. R. : "Robust Estimation : A Condensed Partial Survey", *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 27 (1973), 87-104.
- [11] ——— : "The Influence Curve and its Role in Robust Estimation," *Journal of the American Statistical Association*, 69 (1974), 383-394.
- [12] Hoaglin, D. C. and R. E. Welsch : "The Hat Matrix in Regression and ANOVA," *The American Statistician*, 32(1978), 17-22.
- [13] Holland, P. W. and R. E. Welsch : "Robust Regression Using Iteratively Reweighted Least Squares," *Communications in Statistics*, A6 (1977), 813-827.
- [14] Huber, P. J. : "Robust Statistics : A Review," *The Annals of Mathematical Statistics*, 43 (1972), 1041-1067.
- [15] ——— : "Robust Regression : Asymptotics, Conjectures and Monte Carlo," *The Annals of Statistics*, 1 (1973), 799-821.
- [16] ——— : "Robustness and Designs," in *A Survey of Statistical Design and Linear Models* : J. N. Srivastava, ed., Amsterdam : North Holland, 1975, 287-301.
- [17] Krasker, W. S. : "Estimation in Linear Regression Models with Disparate Data Points," *Econometrica*, 48 (1980), 1333-1346.
- [18] Krasker, W. S. and R. E. Welsch : "Efficient Bounded-Influence Regression Estimation," *Journal of the American Statistical Association*, 77 (1982), 595-604.
- [19] Rao, C. R. : *Linear Statistical Inference and Its Applications*, 2nd ed.. New York : John Wiley & Sons, 1973.
- [20] Theil, H. : *Principles in Econometrics*. Amsterdam : North Holland, 1979.